



Paper

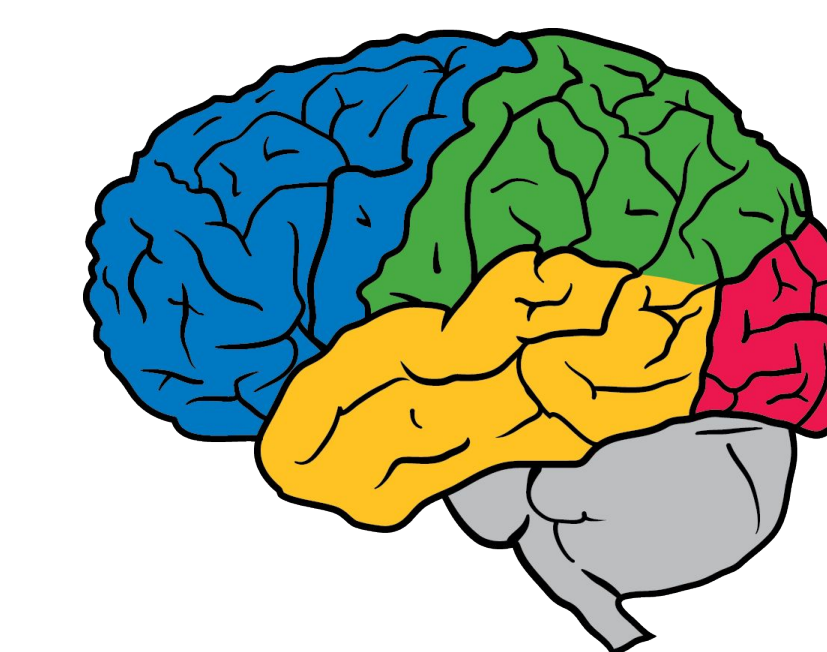


Code

# Wide Neural Networks of Any Depth Evolve as Linear Models Under Gradient Descent

Jaehoon Lee\*, Lechao Xiao\*, Samuel S. Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, Jeffrey Pennington

Google Research, \*Equal contribution, work done as part of the Google AI Residency, g.co/airesidency



NeurIPS 2019; arXiv: 1902.06720

## Are Training Dynamics Tractable?

A longstanding goal in deep learning research has been to precisely characterize **training** and **generalization**. However, the often complex loss landscapes of neural networks have made a theory of learning dynamics elusive. Nevertheless, in the **large width** limit, neural networks evolve as linear models with solvable dynamics.

- Parameter space dynamics:** wide network training dynamics in parameter space are equivalent to the training dynamics of a model which is affine in the collection of all network parameters.
- Sufficient conditions for linearization:** there exists a threshold learning rate  $\eta_{\text{critical}}$  such that gradient descent training of neural networks with learning rate smaller than that threshold are well approximated by their linearization for large width.
- Output distribution dynamics:** the predictions of a neural network throughout gradient descent training converge weakly to a GP as the width goes to infinity. We derive time-dependent expressions for the evolution of this GP and note the differences from the Bayesian posterior GP.
- More in the paper:**
  - Parameterization independence for linearization
  - Momentum, non-square losses (e.g. cross-entropy)
  - Analytic expressions for NTK for Erf and Relu
  - Colab tutorial

## Motivation: Gradient descent learning dynamics of deep neural networks are intractable

- Consider (convex) loss function of neural networks  $\mathcal{L} = \sum_{(x,y) \in \mathcal{D}} \ell(f_t(x, \theta), y)$ .

- Gradient descent (flow) dynamics of the parameters and outputs (logits)

$$\begin{aligned} \dot{\theta}_t &= -\eta \nabla_{\theta} f_t(\mathcal{X})^T \nabla_{f_t(\mathcal{X})} \mathcal{L} \\ \dot{f}_t(\mathcal{X}) &= \nabla_{\theta} f_t(\mathcal{X}) \dot{\theta}_t = -\eta \hat{\Theta}_t(\mathcal{X}, \mathcal{X}) \nabla_{f_t(\mathcal{X})} \mathcal{L} \end{aligned}$$

- With a time evolving **tangent kernel**  $\hat{\Theta}_t = \nabla_{\theta} f_t(\mathcal{X}) \nabla_{\theta} f_t(\mathcal{X})^T$

## Dynamics of linear models are tractable

- For linear models, training dynamics are tractable. Closed form solution for MSE loss:

$$\omega_t = -\nabla_{\theta} f_0(\mathcal{X})^T \hat{\Theta}_0^{-1} (I - e^{-\eta \hat{\Theta}_0 t}) (f_0(\mathcal{X}) - \mathcal{Y}) \quad \text{weight dynamics}$$

$$f_t^{\text{lin}}(\mathcal{X}) = (I - e^{-\eta \hat{\Theta}_0 t}) \mathcal{Y} + e^{-\eta \hat{\Theta}_0 t} f_0(\mathcal{X}) \quad \text{function dynamics}$$

- For some test point  $x$   $f_t^{\text{lin}}(x) = \mu_t(x) + \gamma_t(x)$

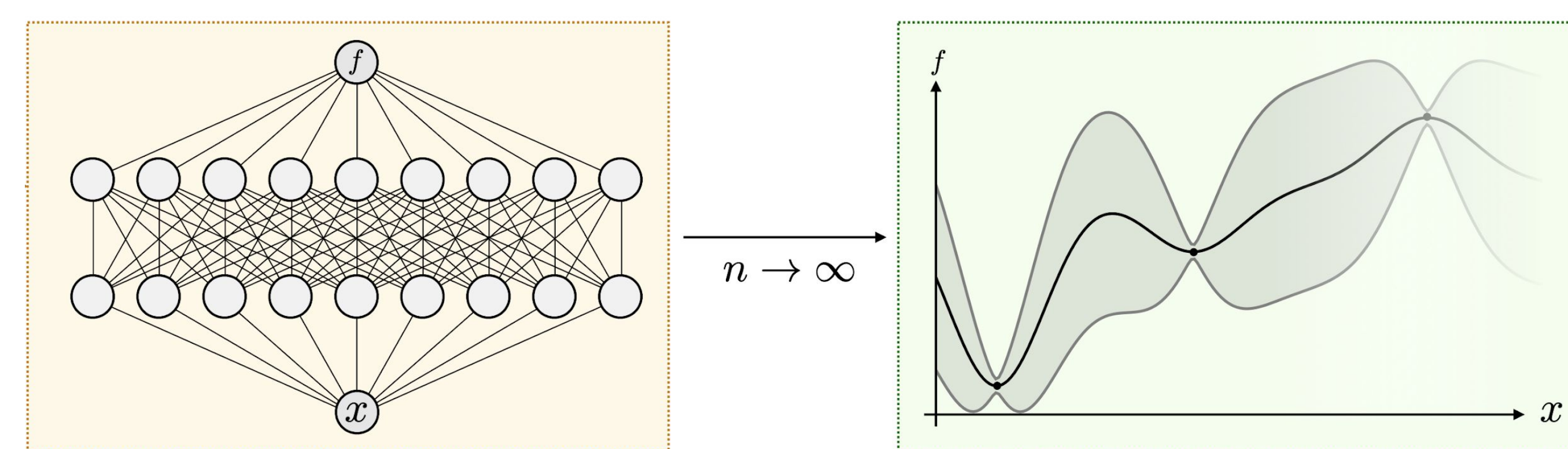
$$\mu_t(x) = \hat{\Theta}_0(x, \mathcal{X}) \hat{\Theta}_0^{-1} (I - e^{-\eta \hat{\Theta}_0 t}) \mathcal{Y} \quad \text{deterministic}$$

$$\gamma_t(x) = f_0(x) - \hat{\Theta}_0(x, \mathcal{X}) \hat{\Theta}_0^{-1} (I - e^{-\eta \hat{\Theta}_0 t}) f_0(\mathcal{X}) \quad \text{stochastic}$$

## Infinite networks are GPs

As the width approaches infinity, the outputs of randomly initialized network converge to a Gaussian Process [1-3]:

$$f_0(\mathcal{X}) \sim \mathcal{N}(0, \mathcal{K}(\mathcal{X}, \mathcal{X})) \quad \mathcal{K}^{i,j}(x, x') = \text{p-lim}_{n \rightarrow \infty} \mathbb{E} [f_0^i(x) f_0^j(x')]$$



This GP is transformed by gradient descent throughout the training process, leading to GP behavior after training.

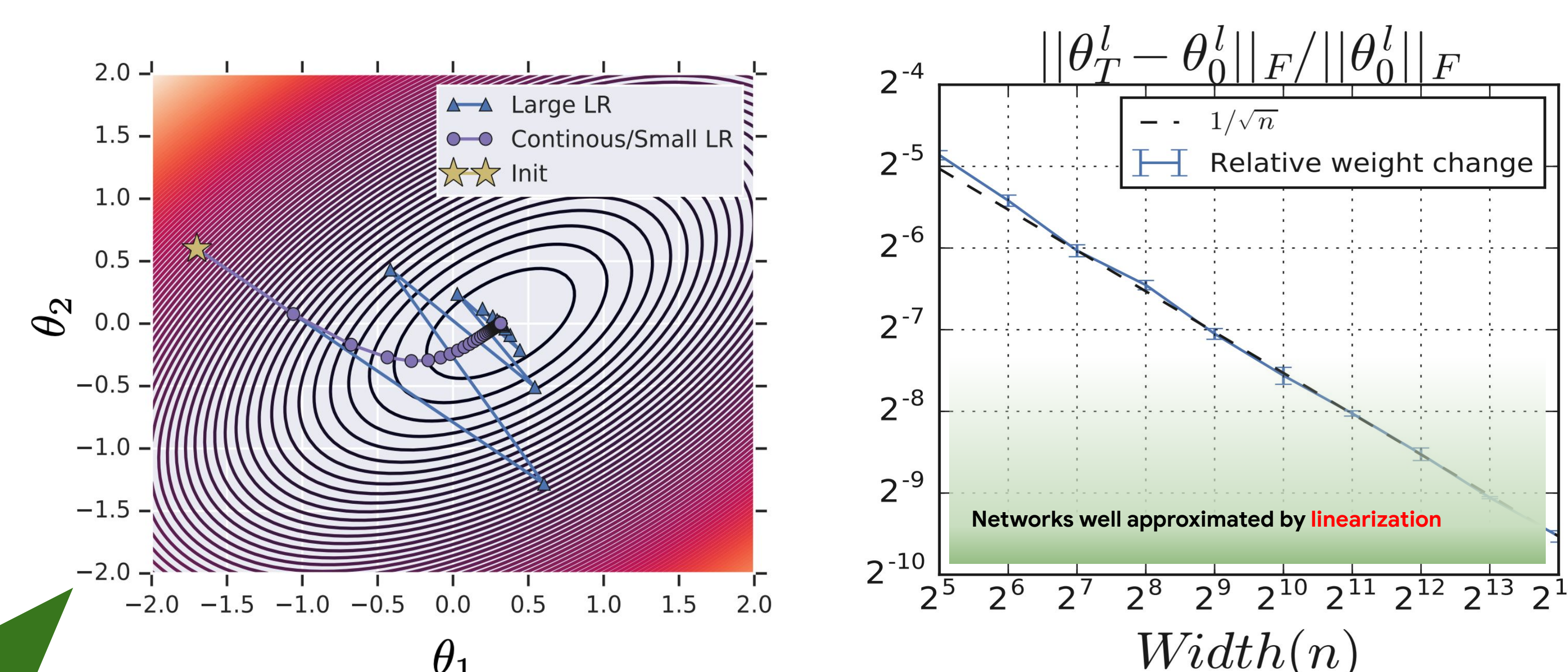
## Infinite networks are linearized models

- When the network is sufficiently **wide**, one can approximate a deep neural network by its first order Taylor expansion (**linearization**) at initialization (t=0)

$$f_t^{\text{lin}}(x) \equiv f_0(x) + \nabla_{\theta} f_0(x) \omega_t \quad \omega_t \equiv \theta_t - \theta_0$$

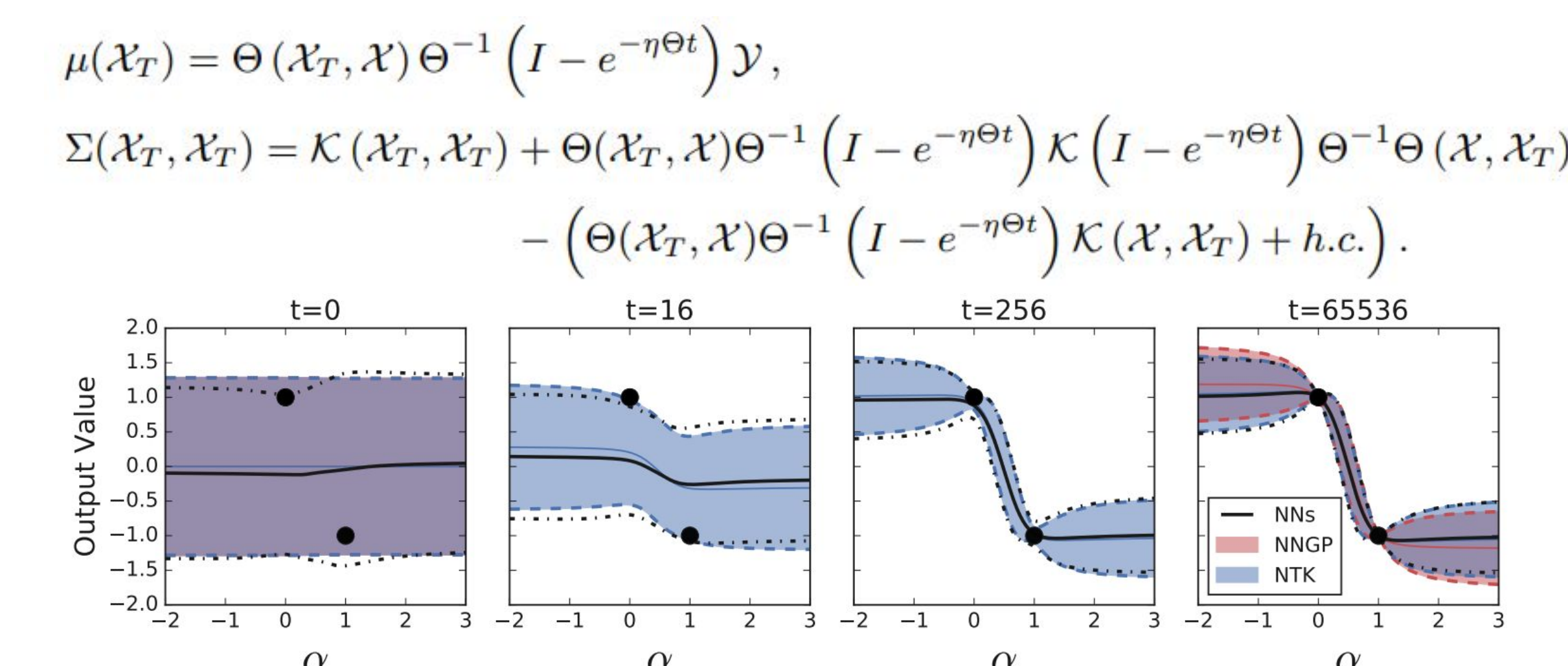
- When learning rate is sufficiently small, linearization becomes more accurate as the width increases:

$$\sup_{t \geq 0} \|f_t(x) - f_t^{\text{lin}}(x)\|_2, \sup_{t \geq 0} \frac{\|\theta_t - \theta_0\|_2}{\sqrt{n}}, \sup_{t \geq 0} \|\hat{\Theta}_t - \hat{\Theta}_0\|_F = \mathcal{O}(n^{-\frac{1}{2}}), \text{ as } n \rightarrow \infty$$

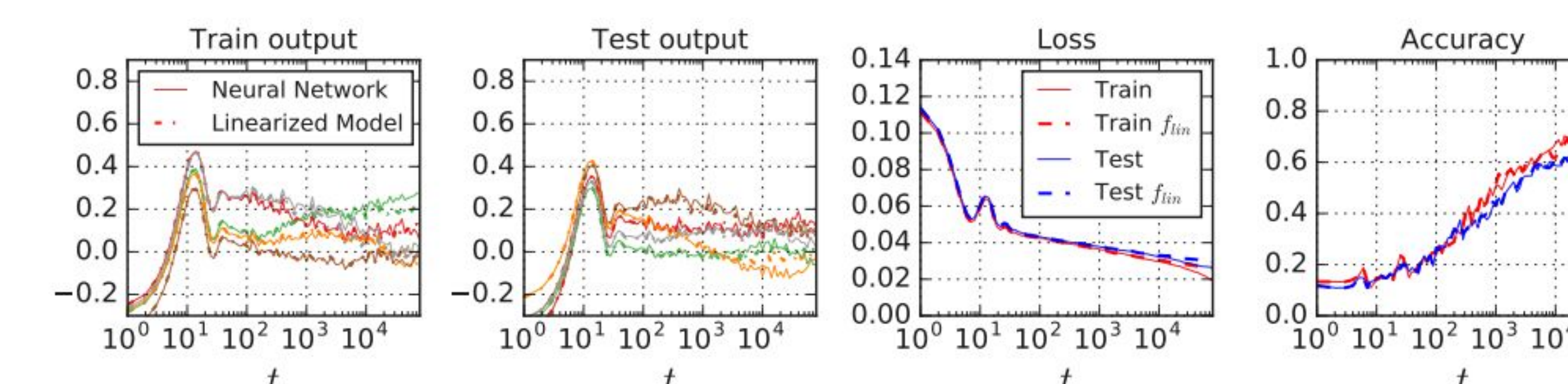


## Experiments

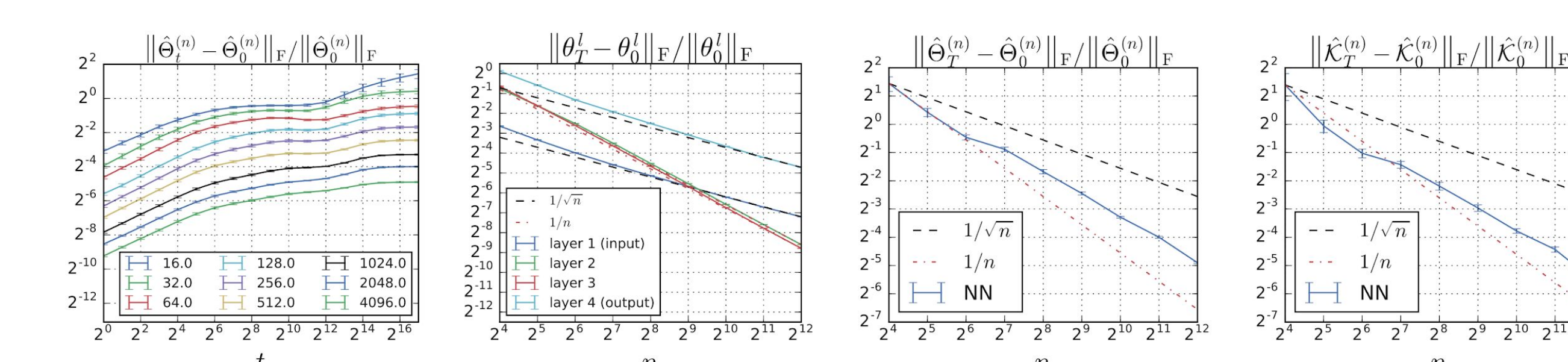
- Gaussian process at initialization leads to Gaussian distribution during training with mean and covariance



- A wide residual network and its linearization behave similarly when both are trained by SGD with momentum on MSE loss on full CIFAR-10



- Relative Frobenius norm change during training



## Relation to other approaches

	Weight Space	Function Space
Bayesian Training	$p(z^*   x^*, \mathcal{D}) = \int d\theta p(z^*   x^*, \theta) p(\theta   \mathcal{D})$ <b>Bayesian Neural Network</b>	$p(z) = \mathcal{GP}(z; 0, K(x))$ $z(x^*)   \mathcal{D} \sim \mathcal{GP}(z_t(x^*)   \mathcal{D})$ <b>NNGP</b>
Gradient Descent	$z(x; \theta_t) = z(x; \theta_0) + \nabla_{\theta} z(x; \theta_0) (\theta_t - \theta_0) + \mathcal{O}(\frac{1}{\sqrt{n}})$ $\partial_{\theta} \theta_t = -\eta (\nabla_{\theta} z_t)^T \nabla_{z_t} \mathcal{L}(z_t(x))$ <b>Linearization</b>	$\partial_t z_t(x) = -\eta \hat{\Theta} \nabla_{z_t} \mathcal{L}(z_t(x))$ $z(x^*)   \mathcal{D} \sim \mathcal{GP}(z_t(x^*) ; \mu_t(\mathcal{D}), \Sigma_t(\mathcal{D}))$ <b>Neural Tangent Kernel</b>

## References

- Neal, Radford M. "Bayesian learning for neural networks", Ph.D. Thesis, University of Toronto, 1994.
- Lee, J. and Bahri, Y., et al. "Deep neural networks as gaussian processes", arXiv:1711.00165, ICLR 2018.
- Matthews, Alexander G. de G., et al. "Gaussian process behaviour in wide deep neural networks", arXiv:1804.11271, ICLR 2018.
- Jacot, A., et al. "Neural Tangent Kernel: Convergence and Generalization in Neural Networks", arXiv:1806.07572, NeurIPS 2018.
- Novak R. and Xiao L., et al. "Neural Tangents: Fast and Easy Infinite Neural Networks in Python", arXiv:1912.02803.

Want to automate NNGP/NTK computation for various architectures?  
: see **"Neural Tangents"** [github.com/google/neural-tangents](https://github.com/google/neural-tangents) [5]!

